

# A Proposal of a Privacy-preserving Questionnaire by Non-deterministic Information and Its Analysis

著者	Tsumoto Shusaku, Nakata Michinori, Sakai Hiroshi, Liu Chenxi
journal or publication title	2016 IEEE International Conference on Big Data (Big Data)
page range	1956-1965
year	2016-12-06
URL	<a href="http://hdl.handle.net/10228/00006129">http://hdl.handle.net/10228/00006129</a>

doi: [info:doi/10.1109/BigData.2016.7840817](https://doi.org/10.1109/BigData.2016.7840817)

# A Proposal of a Privacy-preserving Questionnaire by Non-deterministic Information and Its Analysis

Hiroshi Sakai, Chenxi Liu  
Department of Applied Science,  
Graduate School of Engineering,  
Kyushu Institute of Technology,  
Tobata, Kitakyushu 804-8550, Japan  
E-mail: sakai@mns.kyutech.ac.jp

Michinori Nakata  
Faculty of Management and  
Information Science,  
Josai International University,  
Togane, Chiba 283-8555, Japan  
E-mail: nakatam@ieee.org

Shusaku Tsumoto  
Department of Medical Informatics,  
School of Medicine, Shimane University,  
Enya-cho, Izumo,  
Shimane 693-8501, Japan  
E-mail: tsumoto@med.shimane-u.ac.jp

**Abstract**—We focus on a questionnaire consisting of three-choice question or multiple-choice question, and propose a privacy-preserving questionnaire by non-deterministic information. Each respondent usually answers one choice from the multiple choices, and each choice is stored as a tuple in a table data. The organizer of this questionnaire analyzes the table data set, and obtains rules and the tendency. If this table data set contains personal information, the organizer needs to employ the analytical procedures with the privacy-preserving functionality.

In this paper, we propose a new framework that each respondent intentionally answers non-deterministic information instead of deterministic information. For example, he answers ‘either A, B, or C’ instead of the actual choice A, and he intentionally dilutes his choice. This may be the similar concept on the  $k$ -anonymity. Non-deterministic information will be desirable for preserving each respondent’s information.

We follow the framework of *Rough Non-deterministic Information Analysis* (RNIA), and apply RNIA to the privacy-preserving questionnaire by non-deterministic information. In the current data mining algorithms, the tuples with non-deterministic information may be removed based on the data cleaning process. However, RNIA can handle such tuples as well as the tuples with deterministic information. By using RNIA, we can consider new types of privacy-preserving questionnaire.

## I. INTRODUCTION

We are coping with rough set based information incompleteness, missing values, and data mining [7], [10], [11], [12], [14], [15], [16], [18], [19], [26], [27] in table data sets, and we propose a framework of a privacy-preserving questionnaire in this paper. The idea is simple, namely, each respondent may answer non-deterministic information [14] instead of deterministic information, like ‘either A, B, or C’ instead of the actual one value. Each respondent can dilute his actual answer, and he can preserve his personal information. Of course, such questionnaire will be more privacy-preserved, and we analyze such questionnaire based on *Rough Non-deterministic Information Analysis* (RNIA). Figure 1 shows the total chart for a standard questionnaire, and Figure 2 shows the total chart for the proposing questionnaire.

Recently, the privacy issue on data engineering is very important, and this is often dealt as privacy-preserving data mining [1], [2], [5], [6], [9], [13], [22]. In [1], [2], several ap-

Deterministic Information in Questionnaire

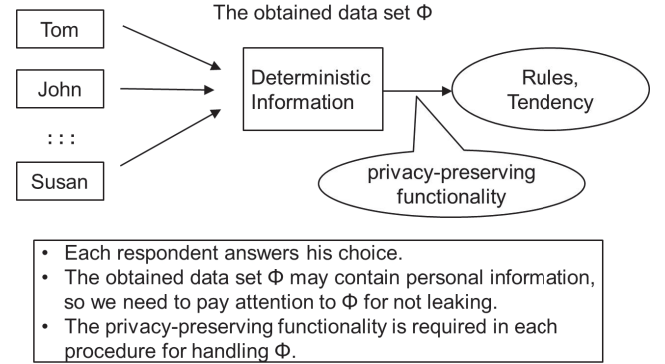


Fig. 1. The total chart handling data set  $\Phi$  with deterministic information.

Non-deterministic Information in Questionnaire

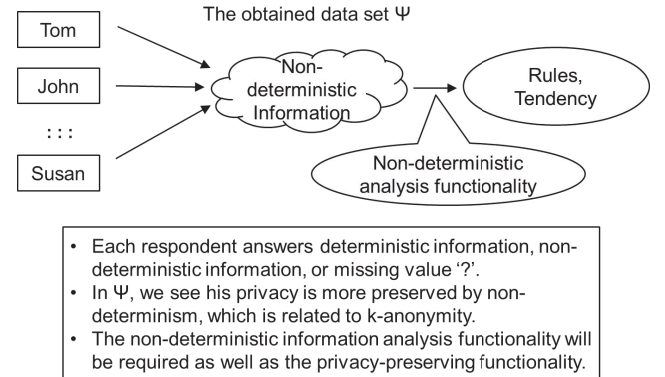


Fig. 2. The total chart handling data set  $\Psi$  with non-deterministic information.

proaches, for example, randomization,  $k$ -anonymization, distributed privacy-preserving data mining, etc. are summarized.

In the randomization method, the noise is added to data for masking the attribute values [5]. This randomization seems to be closely related to the proposing questionnaire. In [5], the

organizer of the questionnaire adds noise to data, but each respondent adds noise to data in our questionnaire. Since the obtained data set by the proposing questionnaire stores vague information, we may have the weakened tendency. However, this will follow the description that the purpose of data mining is to obtain the general tendency of all respondents, and it is not to obtain each personal information [9].

The  $k$ -anonymity model was developed for not to identify any individual records [1], [2]. In our proposal, each respondent may answer ‘either A, B, or C’, which will be corresponding to 3-anonymity in an answer. This will be convenient for each respondent, because we often have information leaks.

In [6], [13], a privacy-preserving web-based questionnaire is investigated, and the issue on a secure protocol for handling the distributed data sets is considered. However, this is the different framework from our proposal.

In our proposal, the main issue is data mining in tables with non-deterministic information, so our proposal will be different from the traditional research on privacy-preserving data mining. However, we think data mining in tables with non-deterministic information is another approach to privacy-preserving.

This paper is organized as follows. In Section 2, we will propose a privacy-preserving questionnaire, and define a questionnaire *QUEST\_Det* and a questionnaire *QUEST\_Non-Det*. In Section 3, we consider NIS-Apriori based rule generation and a prototype system in SQL for handling *QUEST\_Non-Det*. In Section 4, we discuss the merit and the demerit for *QUEST\_Det* and *QUEST\_Non-Det*. In Section 5, we apply the prototype in SQL to Mammographic data set and Lenses data set [8]. In Section 6, we show the use of NIS-Apriori based rule generation by *getRNA* software tool opened in the web [28]. Finally, we conclude the possibility of applying non-deterministic information to a privacy-preserving questionnaire, and clarify the next research.

## II. A PROPOSAL OF A PRIVACY-PRESERVING QUESTIONNAIRE

Let us consider the following questionnaire *QUEST\_Det* by the organizer.

- 1) The questionnaire consists of some questions, which are multiple choices.
- 2) Each respondent answers one choice (deterministic information) from the multiple choices in each question.
- 3) The organizer analyzes the table data set, and obtains the rules for knowing the tendency in the respondents.

### A. A Case of Deterministic Information in a Questionnaire

In this subsection, we consider three-choice questions for simplicity. We often have a questionnaire *QUEST\_Det*, and the organizer may have a table data set  $\Phi$  in Table I. We may call a table data set with definite information a *Deterministic Information System* (DIS) [14], [15], [16], [18], [19]. The organizer will have the public opinion and the tendency by analyzing  $\Phi$ .

TABLE I  
THE OBTAINED DATA SET  $\Phi$  BASED ON THREE-CHOICE QUESTIONS.

<i>Respondent</i>	<i>q1</i>	<i>q2</i>	<i>q3</i>	<i>q4</i>
$r_1$	1	1	1	1
$r_2$	1	2	1	1
$r_3$	2	2	2	2
$r_4$	1	2	2	3
$r_5$	1	3	2	3
$r_6$	2	3	3	3

TABLE II  
THE OBTAINED DATA SETS  $\Psi$ .

<i>Respondent</i>	<i>q1</i>	<i>q2</i>	<i>q3</i>	<i>q4</i>
$r_1$	1	?	1	1
$r_2$	1	2	1	1
$r_3$	2	{1, 2}	2	2
$r_4$	1	2	2	3
$r_5$	{1, 2}	{1, 2, 3}	2	{2, 3}
$r_6$	2	3	3	3

Related to the obtained data set  $\Phi$ , we need to pay attention to the following.

- 1) If  $\Phi$  contains personal and privacy information, it is necessary to assure the security of  $\Phi$ . The privacy-preserving functionality is required in each procedure for handling  $\Phi$ .
- 2) Usually, the security of the obtained data set will be managed by the organizer, however there have been the frequent leaks of data sets with personal information.
- 3) For each respondent, the most convenient answer for preserving his privacy may be the choice of either ‘no answer’ or ‘either A, B, or C’ (non-deterministic information) instead of the actual choice 1.

### B. A Case of Non-deterministic Information in a Questionnaire

In order to preserve his personal information intentionally, we employ non-deterministic information instead of deterministic information. We define each non-deterministic information as a set  $S$  of choices, and we interpret  $S$  as that either an element of  $S$  is the actual choice but we do not know it. For example, non-deterministic information  $\{1, 2\}$  means the actual choice is ‘either 1 or 2’. Non-deterministic information will be the similar concept of the  $k$ -anonymity, and it does not give definite information.

If the organizer of the questionnaire agrees with the use of non-deterministic information, each respondent can intentionally preserve his privacy in the question. Let us suppose we have table  $\Psi$  in Table II. In this case, the answers to  $q_2$  by  $r_1$  and  $r_5$  are semantically the same, because we identify the response ‘no answer’ expressed by ‘?’ symbol with  $\{1, 2, 3\}$  in the three-choice question. We may call a table data set with non-deterministic information a *Non-deterministic Information System* (NIS) [14], [15], [16], [18], [19].

Similarly to  $\Phi$ , we need to pay attention to the following about  $\Psi$ .

- 1) If we have  $\Psi$  instead of  $\Phi$ , personal information will be more preserved.

- 2) However, if we employ the traditional data mining algorithm, the tuples  $r_1$ ,  $r_3$ , and  $r_5$  may be removed from  $\Psi$  by the data cleaning process.
- 3) Therefore, if there are lots of tuples with ? or non-deterministic information in a table data set, the number of considerable tuples may become the small number of tuples. Most of tuples may be ignored by the data cleaning process.
- 4) The research on the data mining algorithm for handling tables like  $\Psi$  will be one solution for preserving personal information.

### C. A Proposal of a Privacy-preserving Questionnaire by Non-deterministic Information

*Proposal 1:* We propose the following questionnaire *QUEST\_Non-Det* by the organizer.

- 1) The questionnaire consists of some questions, which are multiple choices.
- 2) Each respondent may answer non-deterministic information for an inconvenient question. In this case, we may have tables like  $\Psi$  in Table II.
- 3) Each respondent intentionally preserves his personal information by using non-deterministic information.
- 4) *QUEST\_Non-Det* may have vague information than *QUEST\_Det*, so the organizer may not have the precise tendency of the respondents. However, the purpose of the questionnaire seems to know the overview of the respondents' tendency. Therefore, it will be useful to consider *QUEST\_Non-Det* and its analysis.

However, we need new data mining algorithms for analyzing *QUEST\_Non-Det*. In the subsequent section, we consider the analytical method for *QUEST\_Non-Det*.

## III. THE ENVIRONMENT OF DATA ANALYSIS FOR QUESTIONNAIRE *QUEST\_Non-Det*

This section follows the framework of Rough Non-deterministic Information Analysis (RNIA) [19] and apply it to analyzing Questionnaires *QUEST\_Det* and *QUEST\_Non-Det*.

### A. Rules in DIS

This subsection considers rules in DIS, i.e., *QUEST\_Det*. We usually fix an attribute *Dec* as the decision attribute, and handle a pair  $[A, v]$  of the attribute *A* and its attribute value *v*, which we call a *descriptor*. An *implication* is a formula  $\tau : \wedge_i [A_i, v_i] \Rightarrow [Dec, val]$ , and we see an implication satisfying some constraints as a *rule*. In most of work on rule generation, the next two constraints are employed [15], [16], [26], [27], and we also employ them.

For two threshold values  $0 < \alpha, \beta \leq 1.0$ ,  
 $support(\tau) (= Num(\tau) / Number\_of\_the\_tuples) \geq \alpha$ ,  
 $accuracy(\tau) (= Num(\tau) / Num(\wedge_i [A_i, v_i])) \geq \beta$ .

$Num(F)$  is the number of objects supporting a formula *F*. (1)

*Definition 1:* For DIS with a decision attribute *Dec*, and threshold values  $\alpha$  and  $\beta$ , an implication  $\tau$  is a *rule*, if  $\tau$  satisfies both  $support(\tau) \geq \alpha$  and  $accuracy(\tau) \geq \beta$ .

object	attrib	value
r1	q1	1
r1	q2	1
r1	q3	1
r1	q4	1
r2	q1	1
r2	q2	2
r2	q3	1
r2	q4	1
r3	q1	2
r3	q2	2

con1
con2
con3
condi
deci
rdf
rest1
rest2
rest3
rule1
rule2
rule3

Fig. 3. A part of a table  $\Phi_{RDF}$  defined by  $\Phi$ .

Fig. 4. The generated all tables.

Let us consider an implication  $\tau : [q1, 1] \wedge [q3, 2] \Rightarrow [q4, 3]$  in Table I. Since  $\tau$  is supported by the objects  $r_4$  and  $r_5$ , we have  $support(\tau) = 2/6 = 1/3$ , and the formula  $[q1, 1] \wedge [q3, 2]$  is also supported by the objects  $r_4$  and  $r_5$ . Therefore,  $accuracy(\tau) = 2/2 = 1.0$  holds. The  $support(\tau)$  value means the ratio on the occurrence of  $\tau$ , and the  $accuracy(\tau)$  value means the ratio on the consistency of  $\tau$ .

### B. Apriori-based Rule Generation in DIS

This subsection describes rule generation in DIS toward rule generation in NIS.

1) *An Implemented Software in SQL for Handling DISs:* In RNIA, the Apriori algorithm for the transaction data [3], [4], [23] is adjusted to the algorithm for table data sets. Here, each item in the transaction data is identified with a descriptor  $[A, v]$ . We recently employed the environment phpMyAdmin [17], and implemented this Apriori algorithm in SQL procedures based on [19]. By using the actual execution of  $\Phi$  in Table I, we describe the Apriori-based rule generation.

2) *Data Sets in the RDF Format:* We employ the RDF format [24], [25]. Figure 3 is a part of  $\Phi_{RDF}$  in the RDF format of  $\Phi$ . In the usual csv data sets, the attribute value *val* is assigned to the pair of the object *r* and the attribute *attrib*, namely we may see the csv data set is a set of all triplet  $(r, attrib, value)$ .  $\Phi_{RDF}$  is a set of all triplet  $(r, attrib, value)$ , and each tuple in  $\Phi_{RDF}$  corresponds to a descriptor. It is easy to generate a table data set in the RDF format from a csv file.

3) *Rule Generation in Table  $\Phi$ :* We fix the decision attribute *Dec*, the threshold values  $\alpha$  and  $\beta$ , then we execute the SQL query. Figure 5 shows the obtained all rules under the condition *Dec*='q4',  $\alpha=0.3$ , and  $\beta=0.8$ .

This execution took about 1 (sec), and we had all tables in Figure 4. In the current program, it is possible to obtain the rules with less than three conditions. In the first step, tables *con1* and *deci* are generated, then the Cartesian product is generated, and finally the implications satisfying the constraints are stored in a table *rule1*. The implication  $\tau$  satisfying  $support(\tau) \geq \alpha$  and  $accuracy(\tau) < \beta$  is stored in a table *rest1*. In the second step, a table *con2* is generated

att1	val1	decision	support	accuracy	
q2	3	3	0.333	1.000	①
q3	1	1	0.333	1.000	②

att1	val1	att2	val2	decision	support	accuracy	
q1	1	q3	2	3	0.333	1.000	③

Fig. 5. The obtained all rules form  $\Phi_{RDF}$ . In the above table, we have rules  $[q2, 3] \Rightarrow [q4, 3]$ ,  $[q3, 1] \Rightarrow [q4, 1]$ , and  $[q1, 1] \wedge [q3, 2] \Rightarrow [q4, 3]$ .

from *rest1*, and the similar procedure is applied to *con2* and *dec1*. The same procedure is also applied to tables *con3* and *dec1* in the third step.

Like this, it is possible to obtain a set of rules, which show us the tendency in  $\Phi$ .

### C. Rules in NIS

This section follows the framework of RNIA, and considers the rules in NIS, i.e., *QUEST\_Non-Det*. In RNIA, the modal concepts, i.e., the certainty and the possibility are considered, and the certain rules and the possible rules are defined by using all tuples in NIS.

1) *Derived DISs from NIS*: In NIS, we replace each set in a table with a value in the set, and define one DIS, which is named a *derived* DIS from NIS. In  $\Psi$ , we have the 72 derived DISs, and we see an actual DIS is in the 72 derived DISs (Figure 6).

2) *Certain Rules and Possible Rules in NIS*: In RNIA, the following rules are defined based on all derived DISs, namely based on all tuples with non-deterministic information.

**Definition 2:** For NIS with a decision attribute *Dec*, and threshold values  $\alpha$  and  $\beta$ , (1) and (2) are given.

- (1) An implication  $\tau$  is a *certain* rule, if  $\tau$  is a rule in each derived DIS.
- (2) An implication  $\tau$  is a *possible* rule, if  $\tau$  is a rule in at least one derived DIS.

The certain rule is a rule in the unknown actual derived DIS, and it is not influenced by information incompleteness. The possible rule may be a rule in the unknown actual derived DIS, and this rule is related to the possibility. There may be a case that both  $[A, v] \Rightarrow [Dec, val_1]$  and  $[A, v] \Rightarrow [Dec, val_2]$  are the possible rules at the same time.

These two types of rules seem to be the natural extension of rules in DIS, and we will be able to know the tendency of the respondents by using the certain rules and the possible rules.

3) *A Problem for Handling the Certain and the Possible rules*: However, we face with the problem that the number of derived DISs increases exponentially. Even in  $\Psi$ , the number of derived DISs is 72 ( $=2^3 \times 3^2$ ). In Mammographic data set in the UCI machine learning repository [8], the number of derived

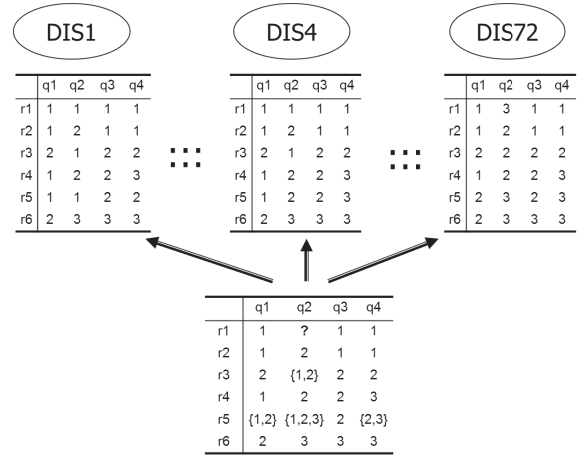


Fig. 6. The 72 derived DISs from  $\Psi$ .

DISs is more than  $10^{100}$ . Therefore, it will be hard to employ the typical method such that we sequentially pick up a derived DIS and examine the constraints.

4) *Theoretical Properties in Rule Generation*: We briefly follow the theoretical properties on rules. Let  $\Omega$  and  $OB$  denote NIS and the set of objects (the set of respondents), respectively. Furthermore, let  $DD(\Omega)$  denote the set of all derived DISs from  $\Omega$ , and the following is defined in  $\Omega$ .

- (1)  $minsupp(\tau) = \min_{\omega \in DD(\Omega)} \{support(\tau) \text{ in } \omega\}$ ,
  - (2)  $minacc(\tau) = \min_{\omega \in DD(\Omega)} \{accuracy(\tau) \text{ in } \omega\}$ ,
  - (3)  $maxsupp(\tau) = \max_{\omega \in DD(\Omega)} \{support(\tau) \text{ in } \omega\}$ ,
  - (4)  $maxacc(\tau) = \max_{\omega \in DD(\Omega)} \{accuracy(\tau) \text{ in } \omega\}$ ,
- If  $\tau$  does not occur in  $\omega$ , we define  $support(\tau) = accuracy(\tau) = 0$  in  $\omega$ .

(2)

The above definitions depend upon each  $\omega \in DD(\Omega)$ , however it is possible to calculate these four values based on [19]. Furthermore, this calculation does not depend upon the number of  $DD(\Omega)$ .

For calculating the above values, at first the following two sets of objects, i.e., *inf* and *sup* blocks, are defined for each descriptor  $[A_i, v_i]$ .

- (1)  $inf([A_i, v_i]) = \{r \in OB \mid \text{the attribute value is deterministic and it is } v_i\}$ ,
- (2)  $inf(\wedge_i [A_i, v_i]) = \cap_i inf([A_i, v_i])$ ,
- (3)  $sup([A_i, v_i]) = \{r \in OB \mid \text{the attribute value is non-deterministic and } v_i \text{ is in the set}\}$ ,
- (4)  $sup(\wedge_i [A_i, v_i]) = \cap_i sup([A_i, v_i])$ .

(3)

Based on the definitions of *inf* and *sup* blocks, the following calculation formula for  $minsupp(\tau)$  to  $maxacc(\tau)$  are given. For NIS  $\Omega$ , let  $\tau$  be an implication  $\wedge_i [A_i, v_i] \Rightarrow [Dec, val]$ . Then, the following is shown in [19].

- (1) For  $\tau$  which occurs in each  $\omega \in DD(\Omega)$ ,  

$$\begin{aligned} \text{minsupp}(\tau) &= |\inf(\wedge_i[A_i, v_i]) \cap \inf([Dec, val])|/|OB|, \\ \text{minacc}(\tau) &= \frac{|\inf(\wedge_i[A_i, v_i]) \cap \inf([Dec, val])|}{|\inf(\wedge_i[A_i, v_i])| + |OUT|}. \end{aligned}$$

Otherwise,  
 $\text{minsupp}(\tau) = \text{minacc}(\tau) = 0.$
- (3)  $\text{maxsupp}(\tau)$   

$$= |\sup(\wedge_i[A_i, v_i]) \cap \sup([Dec, val])|/|OB|.$$
- (4)  $\text{maxacc}(\tau)$   

$$= \frac{|\sup(\wedge_i[A_i, v_i]) \cap \sup([Dec, val])| + |IN|}{|\inf(\wedge_i[A_i, v_i])| + |IN|}.$$
- (5)  $OUT$   

$$= \{\sup(\wedge_i[A_i, v_i]) \setminus \inf(\wedge_i[A_i, v_i])\} \setminus \inf([Dec, val]).$$
- (6)  $IN$   

$$= \{\sup(\wedge_i[A_i, v_i]) \setminus \inf(\wedge_i[A_i, v_i])\} \cap \sup([Dec, val]).$$

Furthermore, the next properties are shown in [19].

- 1) There exists at least one  $\omega_{min} \in DD(\Omega)$  which makes both values of  $\text{support}(\tau)$  and  $\text{accuracy}(\tau)$  the minimum.
- 2) There exists at least one  $\omega_{max} \in DD(\Omega)$  which makes both values of  $\text{support}(\tau)$  and  $\text{accuracy}(\tau)$  the maximum.

Namely, Figure 7 holds for each  $\tau$ , and Theorem 1 is concluded.

**Theorem 1:** The following holds for NIS.

- (1)  $\tau$  is a certain rule, if and only if there is  $\tau$  satisfying  $\text{minsupp}(\tau) \geq \alpha$  and  $\text{minacc}(\tau) \geq \beta$  (Figure 8).
- (2)  $\tau$  is a possible rule, if and only if there is  $\tau$  satisfying  $\text{maxsupp}(\tau) \geq \alpha$  and  $\text{maxacc}(\tau) \geq \beta$  (Figure 9).
- (3) Both conditions in (1) and (2) do not depend upon the number of  $DD(\Omega)$ .

**Example 1:** Let us apply the properties and Theorem 1 to  $\Psi$  under the same condition in Figure 5.

- (1) For  $\tau : [q3, 1] \Rightarrow [q4, 1]$ , the following holds.

$$\begin{aligned} \inf([q3, 1]) &= \sup([q3, 1]) = \{r1, r2\}, \\ \inf([q4, 1]) &= \sup([q4, 1]) = \{r1, r2\}, \\ OUT &= (\sup([q3, 1]) \setminus \inf([q3, 1])) \setminus \inf([q4, 1]) = \emptyset, \\ IN &= (\sup([q3, 1]) \setminus \inf([q3, 1])) \cap \sup([q4, 1]) = \emptyset, \\ \text{minsupp}(\tau) &= |\inf([q3, 1]) \cap \inf([q4, 1])|/|OB| \\ &= |\{r1, r2\}|/6 = 1/3, \\ \text{minacc}(\tau) &= |\{r1, r2\}|/(|\inf([q3, 1])| + |OUT|) \\ &= |\{r1, r2\}|/(|\{r1, r2\}| + |\emptyset|) = 2/2 = 1.0. \end{aligned}$$

Since  $\text{minsupp}(\tau) \geq 0.3$  and  $\text{minacc}(\tau) \geq 0.8$  hold, we conclude the implication  $\tau$  is a certain rule by using Figure 8. Namely, this  $\tau$  is always a rule in each  $\omega \in DD(\Psi)$ . Clearly, the above procedure does not depend upon the number of  $DD(\Psi)$ .

- (2) For  $\tau : [q1, 1] \wedge [q3, 2] \Rightarrow [q4, 3]$ , the following holds.

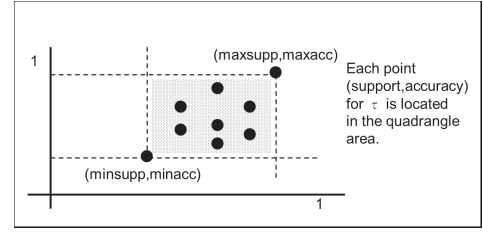


Fig. 7. Each point  $(\text{support}(\tau), \text{accuracy}(\tau))$  by  $\omega$ .

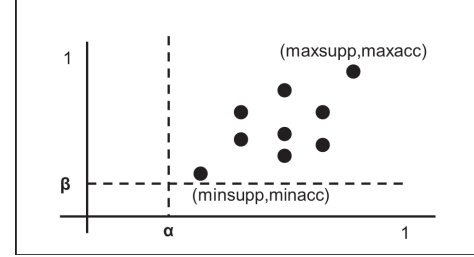


Fig. 8. The case of examining the certain rule.

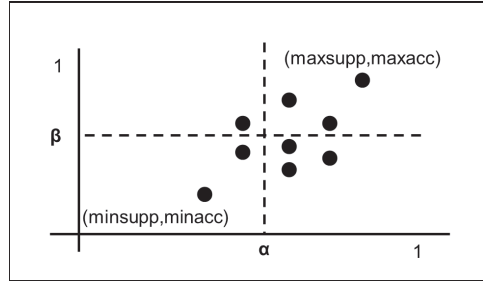


Fig. 9. The case of examining the possible rule.

$$\begin{aligned} \inf([q1, 1]) \cap \inf([q3, 2]) &= \{r1, r2, r4\} \cap \{r3, r4, r5\} = \{r4\}, \\ \sup([q1, 1]) \cap \sup([q3, 2]) &= \{r1, r2, r4, r5\} \cap \{r3, r4, r5\} \\ &= \{r4, r5\}, \\ \inf([q4, 3]) &= \{r4, r6\}, \sup([q4, 3]) = \{r4, r5, r6\}, \\ OUT &= (\{r4, r5\} \setminus \{r4\}) \setminus \{r4, r6\} = \{r5\}, \\ IN &= (\{r5\}) \cap \{r4, r5, r6\} = \{r5\}, \\ \text{minsupp}(\tau) &= |\{r4\} \cap \inf([q4, 3])|/|OB| = |\{r4\}|/6 = 1/6, \\ \text{minacc}(\tau) &= |\{r4\}|/(|\{r4\}| + |\{r5\}|) = 0.5, \\ \text{maxsupp}(\tau) &= |\{r4, r5\} \cap \sup([q4, 3])|/|OB| = |\{r4, r5\}|/6 = 1/3, \\ \text{maxacc}(\tau) &= (|\{r4\} \cap \{r4, r5, r6\}| + |IN|)/(|\{r4\}| + |IN|) \\ &= (|\{r4\}| + |\{r5\}|)/(|\{r4\}| + |\{r5\}|) = 1.0. \end{aligned}$$

Since  $\text{minsupp}(\tau) < 0.3$  and  $\text{minacc}(\tau) < 0.8$  hold, we conclude the implication  $\tau$  is not any certain rule, but  $\text{maxsupp}(\tau) \geq 0.3$  and  $\text{maxacc}(\tau) \geq 0.8$  hold. Therefore, we conclude  $\tau$  is a possible rule.

#### D. Apriori-based Rule Generation in NIS

In RNIA, two rule generation systems are given, namely the certain rule generator and the possible rule generator. Both systems extend the Apriori algorithm in DIS to NIS by using Theorem 1. We describe Apriori-based rule generation by using the actual execution of  $\Psi$ .



### 1) An Overview of NIS-Apriori based Rule Generation:

In the certain rule generator, tables *con1* (the condition part) and *deci* (the decision part) in Figure 4 are generated at first, then for each  $\tau$ ,  $minsupp(\tau)$  and  $minacc(\tau)$  are calculated. By using Theorem 1, the certain rule generator decides  $\tau$  is a certain rule or not. If  $minsupp(\tau) \geq \alpha$  and  $minacc(\tau) \geq \beta$ , this  $\tau$  is stored in a table *crule1*. If  $minsupp(\tau) \geq \alpha$  and  $minacc(\tau) < \beta$ , this  $\tau$  is stored in a table *crest1*. By using *crest1*, the next *ccon2* is generated.

In the possible rule generator,  $maxsupp(\tau)$  and  $maxacc(\tau)$  are calculated by using tables *con1* and *deci*. By using Theorem 1, the possible rule generator decides  $\tau$  is a possible rule or not. If  $maxsupp(\tau) \geq \alpha$  and  $maxacc(\tau) \geq \beta$ , this  $\tau$  is stored in a table *prule1*. If  $maxsupp(\tau) \geq \alpha$  and  $maxacc(\tau) < \beta$ , this  $\tau$  is stored in a table *prest1*. By using *prest1*, a table *pcon2* is generated.

Like this, rule generation in DIS is extended to rule generation in NIS  $\Omega$ . Since every calculation of  $minsupp(\tau)$ ,  $minacc(\tau)$ ,  $maxsupp(\tau)$ , and  $maxacc(\tau)$  does not depend upon  $DD(\Omega)$  and each calculation is the polynomial time order, the extended algorithm will take about the twice complexity of the Apriori algorithm in DIS. In RNIA, this algorithm is called the NIS-Apriori algorithm [19].

### 2) An Implemented Software in SQL for Handling NISs:

We also employed the environment phpMyAdmin [17], and implemented this NIS-Apriori algorithm in SQL procedures based on [19]. By using the actual execution of  $\Psi$ , we describe Apriori-based rule generation in NISs.

3) *Data Sets in the RDF Format:* We employ a table  $\Psi_{NRDF}$  in the NRDF format for  $\Psi$  (Figures 10 and 11). In  $\Psi_{NRDF}$ , each tuple corresponds to a descriptor in one derived DIS. The fourth attribute *det*=1 means the tuple is definite. Otherwise, each tuple comes from non-deterministic information. It is also easy to generate a table data set in the NRDF format.

4) *Rule Generation in Table  $\Psi$ :* We employ the same conditions in  $\Phi$ , namely we fix the condition *Dec*='q4',  $\alpha=0.3$ , and  $\beta=0.8$ . In NIS, we prepared three procedures step1, step2, and step3 in SQL. The step1 generates rules with one condition, the step2 does rules with two conditions, and the step3 does rules with three conditions.

It took about 2 (sec) for executing each step1, step2, and step3. We obtained all tables in Figure 12. Figure 13 and 14 show the obtained all certain rules and all possible rules.

Here, we clarify the relation between rules in DIS and rules in NIS, and we have the following remarks.

**Remark 1:** For NIS  $\Omega$ , let  $\omega$  be a derived DISs from  $\Omega$ . Then, any rule in  $\omega$  is at least a possible rule in  $\Omega$ .

**Remark 2:** The property about *QUEST\_Det* and *QUEST\_Non-Det*

- 1) Some rules in *QUEST\_Det* may be obtained as certain rules in *QUEST\_Non-Det*, for example, ② in Figure 5 and Figure 13.

object	attrib	value	det
r1	q1	1	1
r1	q2	1	3
r1	q2	2	3
r1	q2	3	3
r1	q3	1	1
r1	q4	1	1
r2	q1	1	1
r2	q2	2	1
r2	q3	1	1
r2	q4	1	1

Fig. 10. A part of a table  $\Psi_{NRDF}$  defined by  $\Psi$ .

r3	q1	2	1
r3	q2	1	2
r3	q2	2	2
r3	q3	2	1
r3	q4	2	1
r4	q1	1	1
r4	q2	2	1
r4	q3	2	1
r4	q4	3	1
r5	q1	1	2

Fig. 11. A part of a table  $\Psi_{NRDF}$  defined by  $\Psi$ .

cimpli2
cimpli3
condi
con_des
crest1
crest2
crule1
crule2
crule3
dec_des
imp1
impli1
nrd
pimpli2
pimpli3
prest1
prest2
prule1
prule2
prule3

Fig. 12. The obtained all tables.

att1	val1	deci	deci_value	minsupp	minacc
q3	1	q4	1	0.333	1.000 ②

Fig. 13. The obtained one certain rule  $[q3, 1] \Rightarrow [q4, 1]$ .

att1	val1	deci	deci_value	maxsupp	maxacc	
q2	1	q4	2	0.333	1.000	possible
q2	3	q4	3	0.333	1.000	①
q3	1	q4	1	0.333	1.000	possible

a1	v1	a2	v2	deci	deci_value	maxsupp	maxacc	
q1	1	q3	2	q4	3	0.333	1.000	③
q1	2	q2	2	q4	2	0.333	1.000	possible
q1	2	q3	2	q4	2	0.333	1.000	possible
q2	2	q3	2	q4	3	0.333	1.000	possible

Fig. 14. The obtained all possible rules.

- 2) Any rule in *QUEST\_Det* is at least obtained as a possible rule in *QUEST\_Non-Det*, for example, ① and ③ in Figure 5 and Figure 14.
- 3) Other independent rules may be obtained as possible rules like the rules marked 'possible' in Figure 14.

#### IV. DISCUSSION ABOUT QUEST\_DET AND QUEST\_NON-DET

This section enumerates the merit and the demerit of *QUEST\_Det* and *QUEST\_Non-Det* by using Figure 15. We also discuss about the use of the privacy-preserving questionnaire by non-deterministic information.

*Remark 3: A case of QUEST\_Det  $\Phi$*

Merit:

The precise information is stored in  $\Phi$  ((i) in Figure 15), so it is possible to obtain the correct rules and the tendency of the respondents ((i) and (ii), Figure 15).

Demerit:

Since the precise information is stored in  $\Phi$  ((i) in Figure 15), it is necessary to pay attention to manage  $\Phi$ . Recently, we often have the leaks of data sets. The privacy-preserving functionality is required ((ii) in Figure 15).

*Remark 4: A case of QUEST\_Non-Det  $\Omega$*

Merit:

Information is diluted from  $\Phi$  to  $\Omega$  ((iii) in Figure 15), which defines  $DD(\Omega)$ . This  $\Omega$  is more privacy-preserved than  $\Phi$ . This  $\Omega$  may be effective for the leaks of the data sets. Instead of the rules in  $\Phi$ , the certain rules and the possible rules are obtained in  $\Omega$ . These two rules are defined by all tuples with non-deterministic information, and we have the weakened results from  $\Phi$ , but the privacy in  $\Omega$  is more preserved than that in  $\Phi$ .

Demerit:

The certain rule is a rule in  $\Phi$ , but there may be several possible rules independent from  $\Phi$ . In RNIA, we cannot discriminate the possible rules in  $\Phi$  with the possible rules independent from  $\Phi$ .

*Remark 5: The trade-off relation between  $\Phi$  and  $\Omega$*

- 1) We have the correct results from the precise information, i.e.,  $\Phi$ , and we have the weakened results from the vague information, i.e.,  $\Omega$ .
- 2) On the other hand, information in  $\Omega$  is more privacy-preserved than that in  $\Phi$ .
- 3) Probably, this types of questionnaire has not been considered, because many tuples may be ignored based on the information cleaning process. However, the framework of RNIA gives a new possibility for the privacy-preserving questionnaire.
- 4) In (iv), Figure 15, there seems less software tools except the software in RNIA.

In this subsection, we clarified the merit and the demerit of *QUEST\_Det* and *QUEST\_Non-Det*, and the possibility for the privacy-preserving questionnaire.

#### V. TWO EXAMPLES BY NIS-APRIORI IN SQL

This section gives two examples by the software NIS-Apriori in SQL. The one is data mining from Mammographic

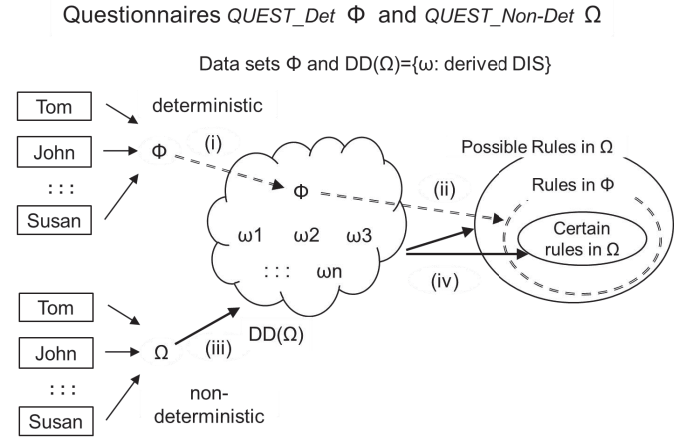


Fig. 15. The survey of *QUEST\_Det* and *QUEST\_Non-Det*.

data set in the UCI machine learning repository [8], and the other is the questionnaire based on Lenses data set in this repository.

##### A. A Case of Mammographic Data Set

Mammographic data set, which handles the cancer data, consists of 960 objects, 6 attributes (*assessment*, *age*, *shape*, *margin*, *density*, and *severity*). We can see the *shape* and *density* as four-choice questions, the *margin* as five-choice question, and the decision *severity* as two-choice (1:benign and 2:malignant) question.

There are about 150 missing values, and we replace each missing value to a set of all choices, and we obtained NIS  $\Sigma$ . If we employ *QUEST\_Non-Det*, we will have the similar data set as  $\Sigma$ . In this case, the number of derived DISs is more than  $10^{100}$ , so it will be hard to obtain the certain rules and the possible rules without Theorem 1. Figure 16 and 17 show the obtained rules from  $\Sigma$ . It took about 20(sec) for *step1* and 7(sec) for *step2*, and we obtained rules with only one condition. Since ④ to ⑦ are certain rules, we conclude that they also hold in the unknown actual derived DIS.

##### B. A Case of the Lenses Data Set

Lenses data set, which handles the contact lenses data, consists of 24 objects, five attributes *age* with three choices 1, 2, 3, *spec*, *asti*, *tear* with two choices 1, 2, and the decision attribute *dec* with three choices 1: hard contact lenses, 2: soft contact lenses, 3: no lenses.

In this data set, there is no missing values, and we see this data set as one DIS  $\pi$ , namely we see  $\pi$  is the obtained data set by *QUEST\_Det*. We randomly added non-deterministic information to  $\pi$ , and generated one NIS  $\Pi$ , which we see the obtained data set by *QUEST\_Non-Det*. In  $\Pi$ , the 30 attribute values in  $\pi$  are changed to non-deterministic information. This means 25% ( $=30/120$ ) information of  $\pi$  is hidden, and  $DD(\Pi)$  consists of 8153726976 ( $=2^{25} \times 3^5$ ) derived DISs. The DIS  $\pi$  is an element of  $DD(\Pi)$ . Figure 18 shows the obtained rules in  $\pi$ , and Figure 19 shows the obtained all possible rules in



att1	val1	deci	deci_value	minsupp	minacc	
assess	4	severity	0	0.445	0.763	④
assess	5	severity	1	0.317	0.869	⑤
margin	1	severity	0	0.329	0.859	⑥
shape	4	severity	1	0.328	0.752	⑦

Fig. 16. The obtained all certain rules from  $\Sigma$  under the condition  $Dec=severity$ ,  $\alpha=0.2$ , and  $\beta=0.7$ .

att1	val1	deci	deci_value	maxsupp	maxacc	
assess	4	severity	0	0.451	0.783	④
assess	5	severity	1	0.330	0.888	⑤
margin	1	severity	0	0.368	0.896	⑥
shape	1	severity	0	0.214	0.844	possible
shape	2	severity	0	0.203	0.848	possible
shape	4	severity	1	0.341	0.794	⑦

Fig. 17. The obtained all possible rules from  $\Sigma$  under the condition  $Dec=severity$ ,  $\alpha=0.2$ , and  $\beta=0.7$ .

att1	val1	decision	support	accuracy	
astigmatic	2	3	0.333	0.667	⑧
spec	2	3	0.333	0.667	⑨
tear	1	3	0.500	1.000	⑩

Fig. 18. The obtained all rules from  $\pi$  under the condition  $Dec=decision$ ,  $\alpha=0.3$ , and  $\beta=0.6$ .

att1	val1	deci	deci_value	maxsupp	maxacc	
age	1	decision	3	0.333	0.667	
age	2	decision	3	0.333	0.727	
astigmatic	1	decision	3	0.458	0.688	
astigmatic	2	decision	3	0.333	0.667	⑧
spec	1	decision	3	0.458	0.688	
spec	2	decision	3	0.458	0.733	⑨
tear	1	decision	3	0.542	1.000	⑩
tear	2	decision	3	0.583	0.609	

Fig. 19. The obtained all possible rules from  $\Pi$  under the condition  $Dec=decision$ ,  $\alpha=0.3$ , and  $\beta=0.6$ .

II. It took about 1 (sec) for generating rules in  $\pi$ , and it took about 2 (sec) in *step1* for  $\Pi$ . We did not obtain any certain rules.

Like this, we will be able to consider rules not only in *QUEST\_Det* but also in *QUEST\_Non-Det*.

## VI. SOME EXAMPLES BY *getRNA* SYSTEM

We have also implemented a software tool *getRNA* (<http://getrnia.org>) [28] in Python based on NIS-Apriori algorithm. This is a demonstrative software tool for showing rule generation in NIS, and we can easily execute some demo files, for example, Soybean.csv, Congress.csv, Cancer.csv, Mammographic.csv, Hepatitis.csv, etc. in Figure

Archived Files: Default ▾		
soybean(L).csv	congress.csv	cancer.csv
Missing: 712 Cases: 5,910 * 10^285 307 * 36 2014-01-09	Missing: 392 Cases: 1,009 * 10^118 435 * 17 2014-01-09	Missing: 0 Cases: 1 699 * 10 2013-10-18
hepatitis.pl	mammo.pl	diluhrs2.pl
Missing: 168 Cases: 2,211 * 10^124 155 * 20 2013-10-03	Missing: 179 Cases: 1,784 * 10^116 980 * 6 2013-09-02	Missing: 13 Cases: 4,147 * 10^04 8 * 4 2013-10-06
salaryNIS.pl	salaryDIS.pl	diluhrs.pl
Missing: 4 Cases: 16	Missing: 0 Cases: 1	Missing: 13 Cases: 2,765 *

Fig. 20. Demo files in *getRNA*.

20. They are all picked up from the UCI machine learning repository.

Figure 21 shows the revised Congress.csv data set, which consists of the 435 objects, the 17 attributes, the 392 missing values, and the number of derived DISs is more than  $10^{100}$ . In the first attribute, the set of the attribute values is  $\{democrat, republican\}$ , and the other set of the attribute values is  $\{yes, no\}$ . Each missing value “?” is changed to a set of all possible values, i.e.,  $\{yes, no\}$ .

Figure 22 shows the screen shot on rule generation, where the constraints are  $support(\tau) \geq 0.5$  and  $accuracy(\tau) \geq 0.7$ . The blue circle implies the minimum point ( $minsupp(\tau), minacc(\tau)$ ), and the orange circle implies the maximum point ( $maxsupp(\tau), maxacc(\tau)$ ). Even though there are more than  $10^{100}$  derived DISs in Congress data set, each point is located in the rectangle area defined by the minimum point and the maximum point. By checking these two characteristic points, we can obtain rules depending upon more than  $10^{100}$  derived DISs. In *getRNA*, we can easily change the constraints, and we obtain rules defined by new constraints. However, *getRNA* handles rules in the form of either  $Con_1 \Rightarrow Decision$  or  $Con_1 \wedge Con_2 \Rightarrow Decision$  for reducing the execution time.

## VII. CONCLUDING REMARKS

We followed the framework of RNA, and proposed new questionnaire by using non-deterministic information. In this questionnaire, each respondent intentionally preserves his privacy, and we can analyze such data sets. Based on [19], we implemented the prototype of NIS-Apriori in SQL, and applied this software to some data sets like Mammographic data set and Lenses data set. Furthermore, we showed the screen shot

of the execution by *getRNIA*. In our proposal, the background is rule generation in NIS, namely NIS-Apriori based rule generator. Based on the execution logs by the prototype in SQL and *getRNIA*, we think that the background is proved to be robust and stable.

Toward the actual application, we need to consider the details below.

(1) About non-deterministic information and the missing values: The use of non-deterministic information and the missing values will be effective for privacy-preserving. However in the actual questionnaire, we will employ binary choices, i.e., ‘either A or B’, in the multi-choice question. Namely, we will handle the answers with 2-anonymity. Even though we can replace a missing value with non-deterministic information ‘either  $A_1, A_2, \dots$ , or  $A_n$ ’, we should escape from this choice, because such choice causes the more weakened results.

(2) About the software: We are coping with the software tool in SQL, because SQL has the high versatility. We simply simulated the NIS-Apriori algorithm by using the procedures in SQL, and realized the next procedures [20], [21],

- (i)  $apri(Dec, |OB|, \alpha, \beta)$  for RDF format,
- (ii)  $step1(Dec, |OB|, \alpha, \beta)$  for NRDF format,
- (iii)  $step2(Dec, |OB|, \alpha, \beta)$  for NRDF format,
- (iv)  $step3(Dec, |OB|, \alpha, \beta)$  for NRDF format.

In NRDF format, if we consider a case that every  $det=1$ , this corresponds to RDF format. We checked some data sets, and examined that each procedure generated the same tables. In the current implementation, we faithfully simulated the NIS-Apriori algorithm, so we had several temporal tables. The current procedures are just the prototype, and we need to refine the procedures.

(3) About the actual questionnaire data analysis: In this paper, we employed several NISs instead of the actual questionnaire, because the proposing questionnaire defines one NIS. Rule generation in NIS is the main issue in the proposing questionnaire. In the next research, we need to handle the actual questionnaire, and it is necessary to consider the analysis depending upon the property of a questionnaire.

#### ACKNOWLEDGMENT

The first author would be grateful to Prof. Dominik Ślęzak for his guidance to SQL. This work is supported by JSPS (Japan Society for the Promotion of Science) KAKENHI Grant Number 26330277.

#### REFERENCES

- [1] Aggarwal, C., Yu, S. (Eds): Privacy-preserving data mining: models and algorithms. Advances in Database Systems 34, Springer (2008)
- [2] Aggarwal, C., Yu, S.: A general survey of privacy-preserving data mining models and algorithms. Advances in Database Systems 34, Springer (2008)
- [3] Agrawal, R., Srikant, R.: Fast algorithms for mining association rules in large databases. Proc. VLDB'94 Morgan Kaufmann, 487–499 (1994)
- [4] Agrawal, R., Mannila, H., Srikant, R., Toivonen, H., Verkamo, A.I.: Fast discovery of association rules. Advances in Knowledge Discovery and Data Mining AAAI/MIT Press, 307–328 (1996)
- [5] Agrawal, R., Srikant, R.: Privacy-Preserving Data Mining. Proc. ACM SIGMOD Conference, 439–450 (2000)

- [6] Bogdanov, D., Talviste, R.: A prototype of online privacy-preserving questionnaire system. Technical report, University of Tartu (2010)
- [7] Clark, P., Grzymala-Busse, J.: An analysis of probabilistic approximations for rule induction from incomplete data sets. Fundamenta Informaticae 132(3), 365–379 (2014)
- [8] Frank, A., Asuncion, A.: UCI machine learning repository. Irvine, CA: University of California, School of Information and Computer Science (2010)  
<http://mllearn.ics.uci.edu/MLRepository.html>
- [9] Kikuchi, H.: Privacy preserving data mining. Proc. FIT2004, (in Japanese) (2004)
- [10] Lipski, W.: On semantic issues connected with incomplete information databases. ACM Transactions on Database Systems 4(3), 262–296 (1979)
- [11] Lipski, W.: On databases with incomplete information. Journal of the ACM 28(1), 41–70 (1981)
- [12] Nakata, M., Sakai, H.: Twofold rough approximations under incomplete information. Int'l. J. General Systems 42(6), 546–571 (2013)
- [13] Nakazato, J., Fujimoto, K., Kikuchi, H.: Privacy preserving web-based questionnaire. Proc. IEEE 19th Int'l. Conf. on Advanced Information Networking and Applications, 285–288 (2005)
- [14] Orłowska, E., Pawlak, Z.: Representation of nondeterministic information. Theoretical Computer Science 29(1-2), 27–39 (1984)
- [15] Pawlak, Z.: Systemy Informacyjne: Podstawy Teoretyczne (in Polish) WNT (1983)
- [16] Pawlak, Z.: Rough Sets: Theoretical aspects of reasoning about data. Kluwer Academic Publishers (1991)
- [17] phpMyAdmin Web Page <http://www.phpmyadmin.net/> (2016)
- [18] Sakai, H., Ishibashi, R., Koba, K., Nakata, M.: Rules and apriori algorithm in non-deterministic information systems. Transactions on Rough Sets 9, 328–350 (2008)
- [19] Sakai, H., Wu, M., Nakata, M.: Apriori-based rule generation in incomplete information databases and non-deterministic information systems. Fundamenta Informaticae 130(3), 343–376 (2014)
- [20] Sakai, H., Liu, C., Zhu X., Nakata, M.: On NIS-Apriori based data mining in SQL. Proc. Int'l. Conf. on Rough Sets 2016, Springer LNAI 9920, 514–524 (2016)
- [21] Sakai, H.: Software Tools for RNIA (Rough Non-deterministic Information Analysis) Web Page (2016)  
<http://www.mns.kyutech.ac.jp/~sakai/RNIA/>
- [22] Sakuma, J., Kobayashi, S.: Privacy-preserving data mining. Journal of the Japanese Society of AI 26(5), 1–11 (2011)
- [23] Sarawagi, S., Thomas, S., Agrawal, R.: Integrating association rule mining with relational database systems: alternatives and implications. Data Mining and Knowledge Discovery 4(2), 89–125 (2000)
- [24] Ślęzak, D., Sakai, H.: Automatic extraction of decision rules from non-deterministic data systems: Theoretical foundations and SQL-based implementation. DTA2009 Springer CCIS Vol.64, 151–162 (2009)
- [25] Świeboda, W., Nguyen, S.: Rough set methods for large and sparse data in EAV format. Proc. IEEE RIVF 2012, 1–6 (2012)
- [26] Tsumoto, S.: Knowledge discovery in clinical databases and evaluation of discovered knowledge in outpatient clinic. Information Sciences 124(1-4), 125–137 (2000)
- [27] Tsumoto, S.: Automated extraction of hierarchical decision rules from clinical databases using rough set model. Expert Systems with Applications 24, 189–197 (2003)
- [28] Wu, M., Nakata, M., Sakai, H.: An overview of the getRNIA system for non-deterministic data. Procedia Computer Science 22, 615–622 (2013) <http://getrnia.org>

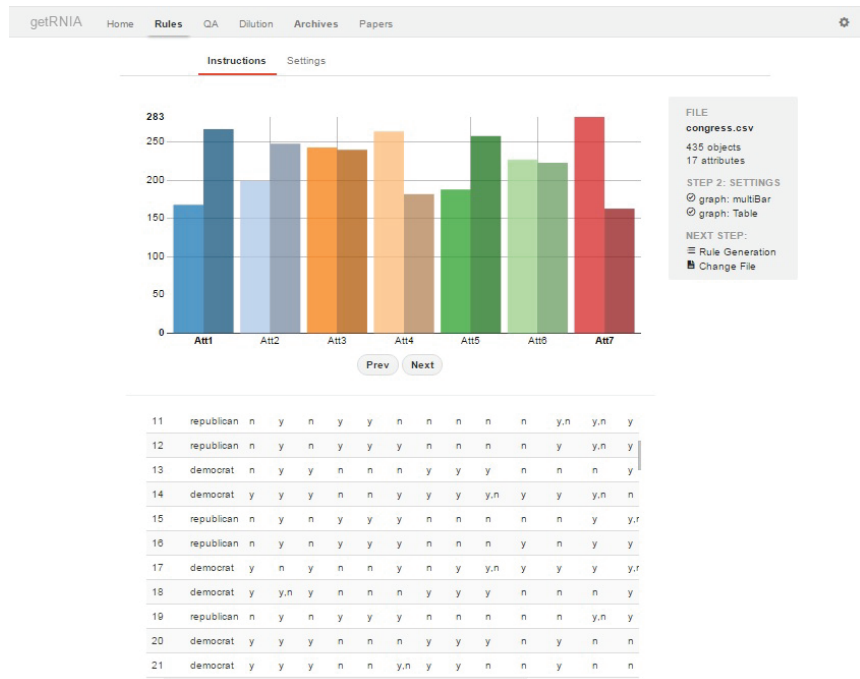


Fig. 21. Congress data set, where each missing value is changed to the possible values {yes,no}. There are more than  $10^{100}$  derived DISs.

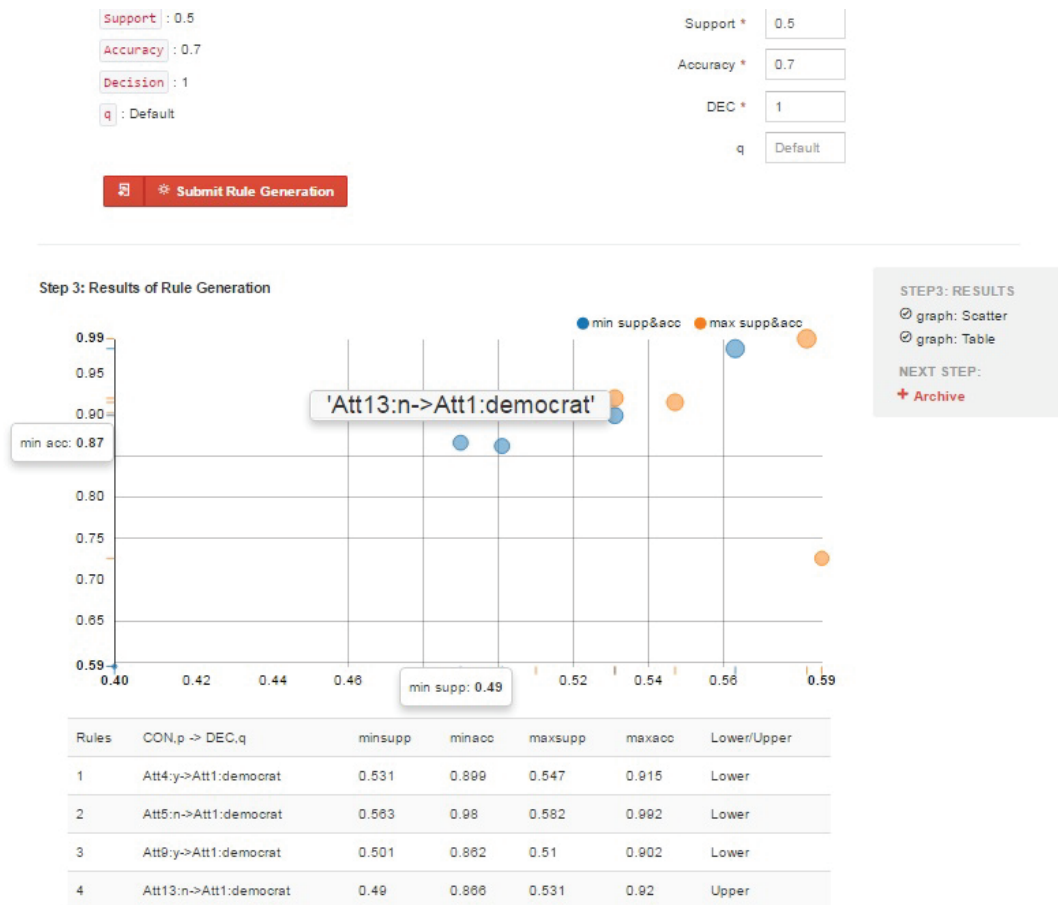


Fig. 22. The generated certain rules (Lower) and the possible rules (Upper). The implication  $[Att13, n] \Rightarrow [Att1, democrat]$  does not satisfy the constraints for a certain rule, but it satisfies the constraints for a possible rule. An implication  $[Att4, y] \Rightarrow [Att1, democrat]$  satisfies the constraints in each of more than  $10^{100}$  derived DISs.